

Metagenomics: Read length matters

K. Eric Wommack¹, Jaysheel Bhavsar¹ and Jacques Ravel^{2*}

¹ Delaware Biotechnology Institute, University of Delaware, 15 Innovation Way, Newark, DE 19711, USA.

² Institute for Genome Sciences, Department of Microbiology & Immunology, University of Maryland School of Medicine, 20 Penn Street, Baltimore, MD 21201, USA.

* To whom correspondence should be addressed. E-mail: jravel@som.umaryland.edu

Cluster of orthologous groups (COG) analysis of short read data simulation experiments (Figs. S1-S3).

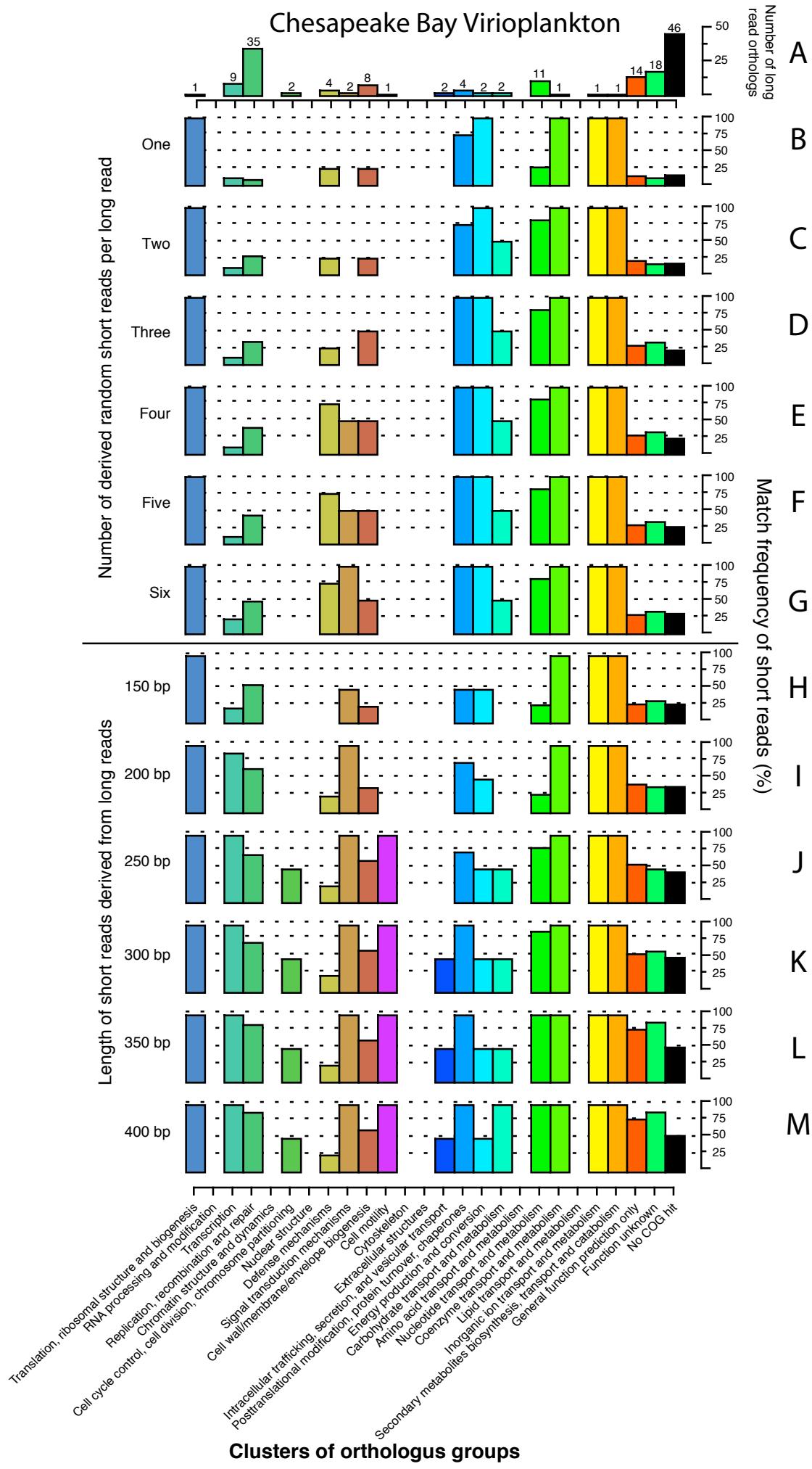


Fig. S1. Distribution of Chesapeake Bay viriplankton BLAST-positive sequences according to clusters of orthologous groups of proteins (COGs). A) Distribution of long read sequences within each COG. Numbers above each bar are the number of long read BLAST-positive sequences within the COG. B-G) Increasing number of 100 bp random derived short read sequences per long read. H-M) Increasing length of single random derived short reads. BLAST-positive short derived sequences expressed as a percent of the long read sequences within each COG. The frequency of long read BLAST-positive sequences that did not belong to a COG was 56% (209 sequences).

Acid Mine Drainage

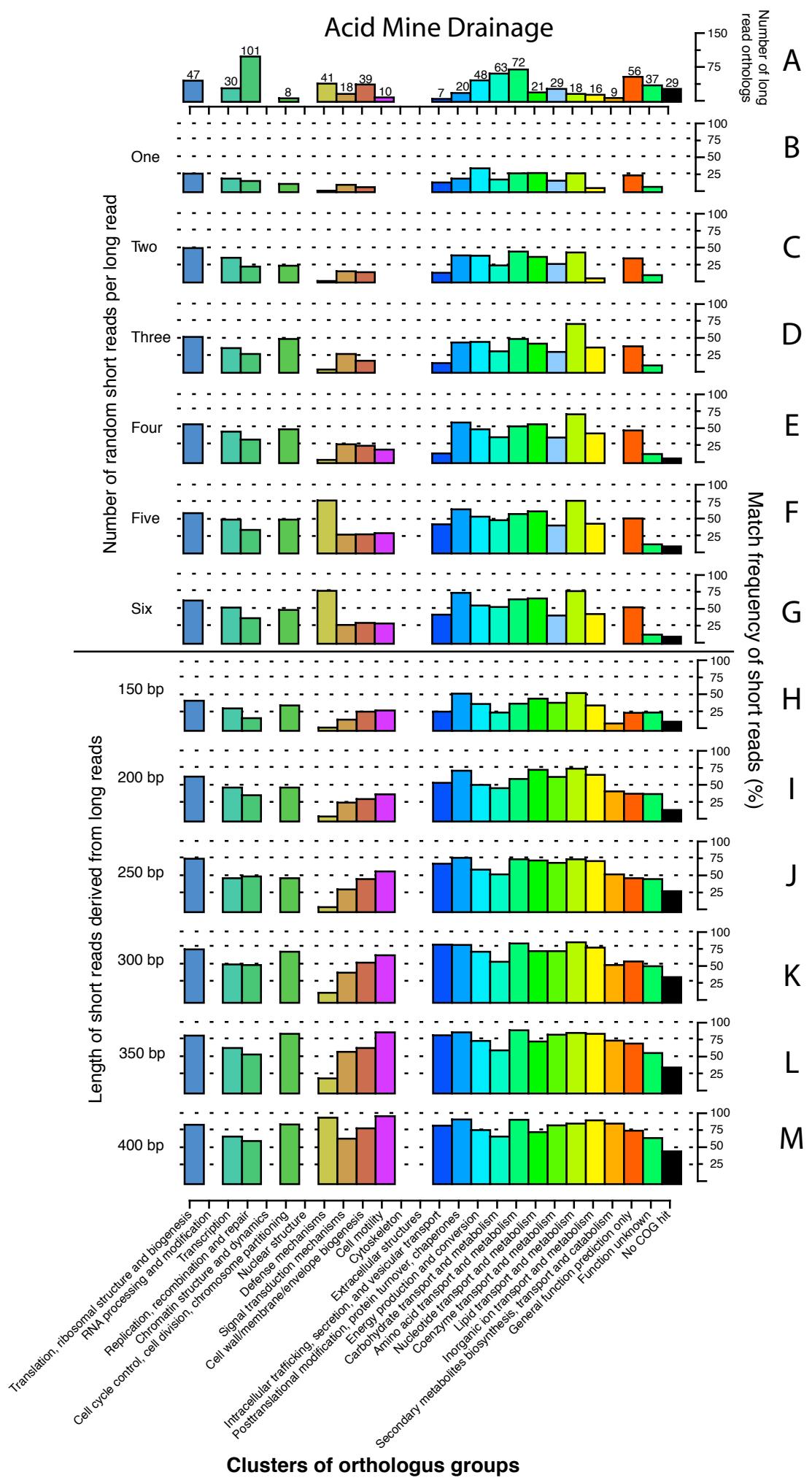


Fig. S2. Distribution of Acid mine drainage microbial metagenome BLAST-positive sequences according to clusters of orthologous proteins (COGs). A) Distribution of long read sequences within each COG. Numbers above each bar are the number of long read BLAST-positive sequences within the COG. B-G) Increasing number of 100 bp random derived short read sequences per long read. H-M) Increasing length of single random derived short reads. BLAST-positive short derived sequences expressed as a percent of the long read sequences within each COG. The frequency of long read BLAST-positive sequences that did not belong to a COG was 13% (107 sequences).

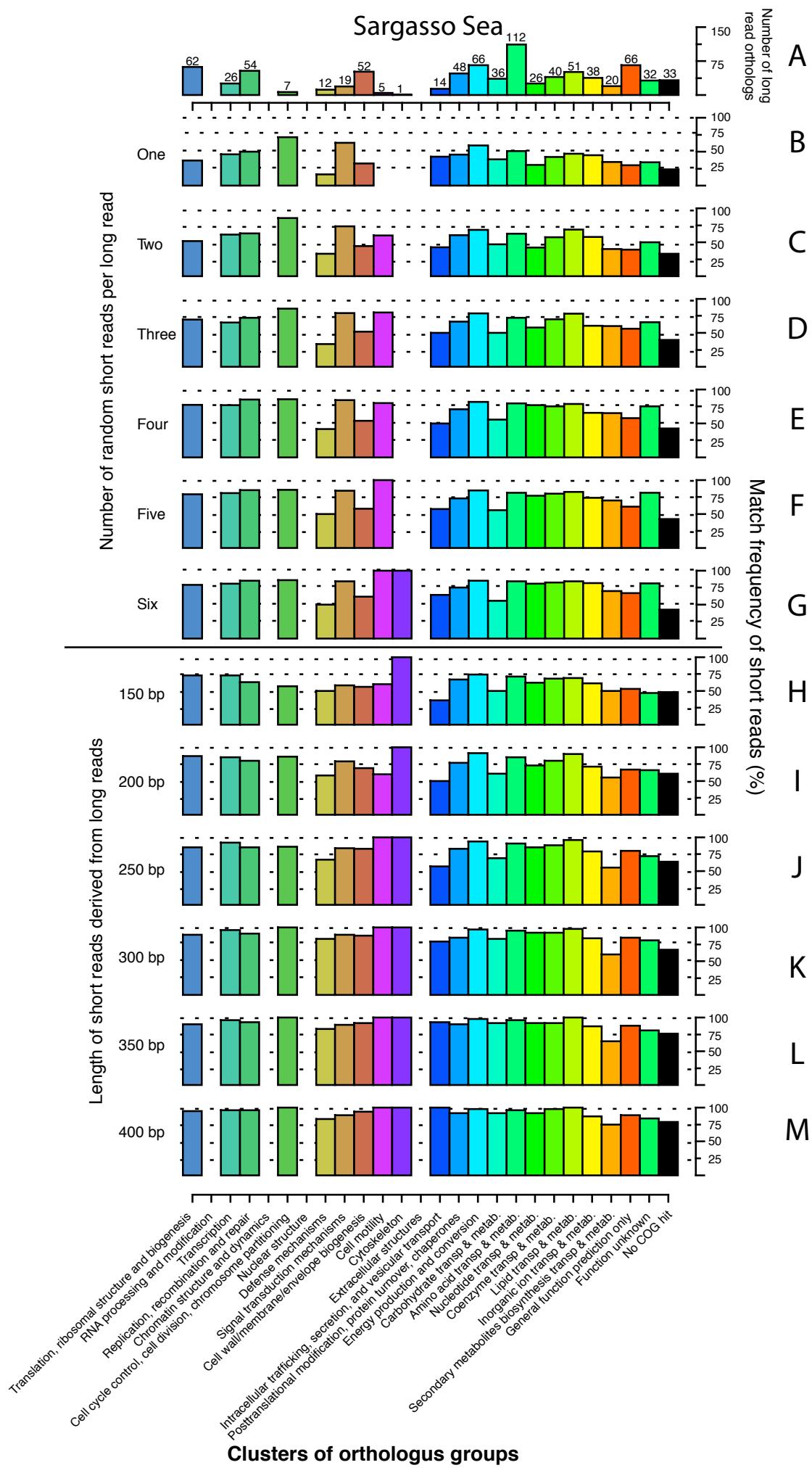


Fig. S3. Distribution of Sargasso Sea microbial metagenome BLAST-positive sequences according to clusters of orthologous proteins (COGs). A) Distribution of long read sequences within each COG. Numbers above each bar are the number of long read BLAST-positive sequences within the COG. B-G) Increasing number of 100 bp random derived short read sequences per long read. H-M) Increasing length of single random derived short reads. BLAST-positive short derived sequences expressed as a percent of the long read sequences within each COG. The frequency of long read BLAST-positive sequences that did not belong to a COG was 5.4% (47 sequences).